

Indirect Reciprocity with Trinary Reputations

Shoma Tanabe¹, Hideyuki Suzuki², and Naoki Masuda^{1,*}

¹ Department of Mathematical Informatics,
The University of Tokyo,
7-3-1 Hongo, Bunkyo, Tokyo 113-8656, Japan

² Institute of Industrial Science,
The University of Tokyo,
4-6-1 Komaba, Meguro, Tokyo 153-8505, Japan

* Corresponding author (masuda@mist.i.u-tokyo.ac.jp)

March 4, 2013

Abstract

Indirect reciprocity is a reputation-based mechanism for cooperation in social dilemma situations when individuals do not repeatedly meet. The conditions under which cooperation based on indirect reciprocity occurs have been examined in great details. Most previous theoretical analysis assumed for mathematical tractability that an individual possesses a binary reputation value, i.e., good or bad, which depends on their past actions and other factors. However, in real situations, reputations of individuals may be multiple valued. Another puzzling discrepancy between the theory and experiments is the status of the so-called image scoring, in which cooperation and defection are judged to be good and bad, respectively, independent of other factors. Such an assessment rule is found in behavioral experiments, whereas it is known to be unstable in theory. In the present study, we fill both gaps by analyzing a trinary reputation model. By an exhaustive search, we identify all the cooperative and stable equilibria composed of a homogeneous population or a heterogeneous population containing two types of players. Some results derived for the trinary reputation model are direct extensions of those for the binary model. However, we find that the trinary model allows cooperation under image scoring under some mild conditions.

1 Introduction

Humans and other animals often cooperate even when cooperation is more costly than defection. In such social dilemma situations, direct reciprocity is among main reasons for cooperation between pairs of individuals that repeatedly meet each other (Trivers, 1971; Axelrod, 1984). However, individuals, in particular humans, cooperate with others even when they seldom meet the same partners more than once, as is the case for large populations. Reputation-based indirect reciprocity (also called downstream reciprocity; we simply call it indirect reciprocity in this paper) seems to be a dominant mechanism that enables cooperation in this situation. In indirect reciprocity, individuals cooperate with others with good reputations and they in turn gain good reputations if they appositely behave (e.g., cooperate) toward somebody else. The conditions under which indirect reciprocity realizes cooperation have been theoretically and numerically clarified in great details (Nowak & Sigmund, 1998a; Nowak & Sigmund, 1998b; Leimar & Hammerstein, 2001; Panchanathan & Boyd, 2003; Mohtashemi & Mui, 2003; Fishman, 2003; Ohtsuki, 2004; Ohtsuki & Iwasa, 2004; Ohtsuki & Iwasa, 2006; Ohtsuki & Iwasa, 2007; Nowak & Sigmund, 2007; Brandt & Sigmund, 2004; Brandt & Sigmund, 2005; Brandt & Sigmund, 2006; Pacheco *et al.*, 2006; Roberts, 2008; Uchida, 2010; Nakamura & Masuda, 2011; Berger, 2011; Sigmund, 2012).

Under the so-called image scoring, cooperation and defection are regarded to be good and bad behavior, respectively (Nowak & Sigmund, 1998a; Nowak & Sigmund, 1998b). Laboratory experiments suggest that humans use image scoring to evaluate others' behavior (Wedekind & Milinski, 2000; Milinski *et al.*, 2001; Seinen & Schram, 2006). However, main theories attain that image scoring does not stabilize cooperation (Leimar & Hammerstein, 2001; Panchanathan & Boyd, 2003; Ohtsuki, 2004; Ohtsuki & Iwasa, 2004; Ohtsuki & Iwasa, 2007; Roberts, 2008). Although some studies have shown the viability of cooperation under image scoring (Nowak & Sigmund, 1998a; Fishman, 2003; Brandt & Sigmund, 2004; Brandt & Sigmund, 2005; Brandt & Sigmund, 2006; Uchida, 2010), the situations in which cooperation occurs are, in our view, quite restricted (but see (Berger, 2011); we discuss this reference in Discussion). Under image scoring, cooperation occurs only when individuals always cooperate in the first round (Nowak & Sigmund, 1998a), unconditional defectors sometimes cooperate (Fishman, 2003), the number of interaction obeys

the binomial distribution (Brandt & Sigmund, 2004) or Poisson distribution (Brandt & Sigmund, 2006), the probability that individuals recognize others' reputations increases in time (Brandt & Sigmund, 2005), or the reputations of individuals are revealed to others with a small probability (Uchida, 2010). Therefore, the reason for the discrepancy between the experiments and theory remains obscure.

For mathematical tractability and possible influences of the first seminal theoretical papers on this subject (Nowak & Sigmund, 1998a; Nowak & Sigmund, 1998b), most theoretical results of indirect reciprocity are derived from the analysis of binary reputation models. In other words, individuals are endowed with the binary reputation, i.e., good (+) or bad (−), depending on the last action toward others and other factors. However, the binary reputation may not be realistic in that only the last behavior of an individual in the social dilemma situation determines the reputation of the individual. In fact, experimental (Wedekind & Milinski, 2000; Milinski *et al.*, 2001; Seinen & Schram, 2006; Milinski *et al.*, 2002; Wedekind & Braithwaite, 2002; Keser, 2003; Bolton & Katok, 2004; Bolton *et al.*, 2005; Engelmann & Fischbacher, 2009) and numerical (Diekmann & Przepiorka, 2005) studies of indirect reciprocity in the context of on-line marketplaces often assume that the reputations are many valued, which complies with the reality of online marketplaces (Resnick & Zeckhauser, 2002; Resnick *et al.*, 2006). More than binary valued reputations have also been employed in numerical studies of indirect reciprocity in theoretical biology literature (Nowak & Sigmund, 1998b; Leimar & Hammerstein, 2001; Mohtashemi & Roberts, 2008). Nevertheless, these studies are not concerned with relationships between the degree of cooperation and the number of the possible reputation values.

In this paper we analyze a trinary reputation model to identify stable populations that realize cooperation. The difference between the present results and those derived from the binary reputation models is remarkable. In particular, we find that image scoring can stabilize cooperation in the trinary reputation model.

2 Model

2.1 Donation Game with Reputations

We consider an infinitely large population. In each generation, the so-called donation game is repeated for sufficiently many rounds. Figure 1(A) illustrates the interaction in each round. Two players are randomly selected from the population, one as donor and the other as recipient, with the equal probability. The donor intends to cooperate (C) or defect (D) toward the recipient according to the action rule σ , which we define below. We assume that the donor misimplements intended C such that the donor actually defects with probability $\epsilon_i > 0$ and that the intended D is always correctly implemented. We seek the possibility of cooperation in the population under this kind of implementation error, which is adverse to cooperation (Panchanathan & Boyd, 2003; Fishman, 2003; Ohtsuki, 2004; Ohtsuki & Iwasa, 2004; Ohtsuki & Iwasa, 2006; Ohtsuki & Iwasa, 2007; Nowak & Sigmund, 2005; Brandt & Sigmund, 2004; Brandt & Sigmund, 2005; Brandt & Sigmund, 2006; Uchida, 2010; Berger, 2011). If the donor implements C, the donor pays cost c , and the recipient obtains benefit b . If the donor implements D, the payoffs to the donor and recipient do not change. We assume that $0 < c < b$ such that the donation game is essentially the prisoner's dilemma.

We assume that each player possesses a reputation that takes one of the three values, i.e., G (Good), N (Neutral), or B (Bad). The action rule σ is a function from the recipient's reputation to the donor's intended action (i.e., C or D). Therefore, there are $2^3 = 8$ action rules, as shown in Fig. 1(B). For example, the AllC and AllD intend C and D regardless of the recipient's reputation, respectively. The so-called generous discriminator (gDisc) intends C when the recipient's reputation is either G or N and D otherwise. The so-called rigorous discriminator (rDisc) intends C when the recipient's reputation is G and D otherwise.

At the end of each round, the observer assesses the donor to be + or - depending on the donor's implemented action (i.e., C or D) and the recipient's reputation (i.e., G, N, or B). Such an assessment rule is a function from the donor's implemented action and the recipient's reputation to the observer's assessment and is called the second-order social norm, which extends the concept in the case of the binary reputation (Leimar & Hammerstein, 2001;

Nowak & Sigmund, 2005). There are $2^{2 \times 3} = 64$ social norms. For example, under the social norm called image scoring (Fig. 1(C)), the observer assigns $+$ and $-$ to the donor's actions C and D, respectively, regardless of the recipient's reputation. Another social norm named scoring-standing is also shown in Fig. 1(C). Finally, we assume that the assessment (i.e., $+$ or $-$) that the donor receives is opposite to the one assigned by the observer with probability $\epsilon_a > 0$ (Ohtsuki & Iwasa, 2004; Ohtsuki & Iwasa, 2007; Nakamura & Masuda, 2011). We assume that the donor's reputation is publicly shared in the population. In other words, the donor's new reputation is instantaneously known to all the players in the population.

We repeat many rounds of the game and define the payoff to each player as the sum of the payoff obtained in the games that involve the player. Because the population is infinite, each pair of players plays the game at most once such that direct reciprocity is excluded.

For a later use, we mention four representative second-order social norms in the binary reputation model (Fig. 1(D)). Image scoring ("scoring" in Fig. 1(D)) does not stabilize cooperation (Panchanathan & Boyd, 2003; Ohtsuki, 2004; Ohtsuki & Iwasa, 2004; Ohtsuki & Iwasa, 2007; Roberts, 2008) unless somewhat strong conditions are met (Nowak & Sigmund, 1998a; Fishman, 2003; Brandt & Sigmund, 2004; Brandt & Sigmund, 2005; Brandt & Sigmund, 2006; Uchida, 2010) (but see (Berger, 2011)). Simple standing ("standing" in Fig. 1(D)) and stern judging ("judging" in Fig. 1(D); also called Kandori (Kandori, 1992)) are known to stabilize cooperation (Ohtsuki & Iwasa, 2004). Shunning stabilizes cooperation under certain conditions (Ohtsuki & Iwasa, 2004; Nakamura & Masuda, 2011).

2.2 Reputation Dynamics

After each round, the donor's reputation is updated on the basis of the observer's assessment. Basically, the donor's reputation shifts upward and downward if the donor receives $+$ and $-$, respectively. When the reputation is binary (i.e., G and B), G and B are equivalent to $+$ and $-$ that the donor receives in the last game, respectively (Nowak & Sigmund, 1998a; Nowak & Sigmund, 2005). However, the relationship between the reputation and assessment is not straightforward in the trinary reputation model. We assume that reputation dynamics obey a Markov chain. We consider the reputation dynamics illustrated in Fig. 2.

In the reputation dynamics shown in Fig. 2(A), which we call the gradual dynamics, the reputation is assumed to move by at most one level in each round. The reputation is unchanged with probability α ($0 \leq \alpha < 1$).

In the reputation dynamics shown in Fig. 2(B), which we call the saltatory dynamics, the donor's reputation can transit from G to B or vice versa in one step. When a G donor receives $-$, the donor's new reputation becomes B and N with probabilities β_d and $1 - \beta_d$, respectively ($0 \leq \beta_d \leq 1$). When a B donor receives $+$, the donor's new reputation becomes G and N with probabilities β_u and $1 - \beta_u$, respectively ($0 \leq \beta_u \leq 1$). When $(\beta_d, \beta_u) = (1, 0)$, the reputation dynamics are similar to the so-called T -period punishment with $T = 2$ (Kandori, 1992) in which players have either state 0 (innocent), 1 (guilty and no repent), ..., or T (guilty and $T - 1$ times of repentant behavior). Guilty players in state 1 regain the innocent state by cooperating with innocent players successive T times. Therefore, the states 0, 1, and 2 are similar to reputations G, B, and N, respectively, in our model. When $(\beta_d, \beta_u) = (0, 1)$, the reputation dynamics are similar to the so-called tolerant scoring (Berger, 2011) because the donor's reputation becomes B if and only if the donor receives $-$ in the last two rounds, and the donor obtains a G reputation if the donor cooperates just once.

3 Analysis Methods

3.1 Homogeneous Populations

We first examine the stability of a homogeneous population of resident players with action rule σ against mutants with different action rules under a given social norm. Let p_G , p_N , and p_B be the probabilities that the reputation of a resident player is G, N, and B, respectively. After a transient of the reputation dynamics, the three probabilities converge to the equilibrium values denoted by p_G^* , p_N^* , and p_B^* . For expository purposes, we focus on the gradual reputation dynamics (Fig. 2(A)) in this section. The following calculations are similar for the saltatory reputation dynamics (Fig. 2(B)); the corresponding results are shown in Appendix A. For the

gradual reputation dynamics, we obtain

$$\begin{cases} p_G^* &= p_G^*[\alpha + (1 - \alpha)\Phi^*] + p_N^*(1 - \alpha)\Phi^*, \\ p_N^* &= p_G^*(1 - \alpha)(1 - \Phi^*) + p_N^*\alpha + p_B^*(1 - \alpha)\Phi^*, \\ p_B^* &= p_N^*(1 - \alpha)(1 - \Phi^*) + p_B^*[\alpha + (1 - \alpha)(1 - \Phi^*)], \end{cases} \quad (1)$$

where Φ^* is the probability that the donor receives + in the equilibrium. Equation (1) and the normalization $p_G^* + p_N^* + p_B^* = 1$ lead to

$$(p_G^*, p_N^*, p_B^*) = \left(\frac{\Phi^{*2}}{1 - \Phi^* + \Phi^{*2}}, \frac{\Phi^*(1 - \Phi^*)}{1 - \Phi^* + \Phi^{*2}}, \frac{(1 - \Phi^*)^2}{1 - \Phi^* + \Phi^{*2}} \right). \quad (2)$$

It should be noted that Eq. (2) and the following results are independent of the value of α . Equation (2) implies that the distribution of the reputation is monotonous in the sense that $p_G^* > p_N^* > p_B^*$ if $\Phi^* > 0.5$ and $p_G^* < p_N^* < p_B^*$ if $\Phi^* < 0.5$.

Φ^* is given by

$$\Phi^* = \sum_{r \in \{G, N, B\}} p_r^* [\zeta_r \Phi_{C,r} + (1 - \zeta_r) \Phi_{D,r}], \quad (3)$$

where ζ_r represents the probability that the donor's implemented action is C when the recipient has reputation $r \in \{G, N, B\}$. ζ_r has a one-to-one correspondence with action rule σ . For example, AllC, gDisc, rDisc, and AllD are equivalent to $(\zeta_G, \zeta_N, \zeta_B) = (1 - \epsilon_i, 1 - \epsilon_i, 1 - \epsilon_i)$, $(1 - \epsilon_i, 1 - \epsilon_i, 0)$, $(1 - \epsilon_i, 0, 0)$, and $(0, 0, 0)$, respectively. $\Phi_{C,r}$ and $\Phi_{D,r}$ in Eq. (3) are the probabilities that the donor receives + when the donor's action is C and D, respectively, and the recipient has reputation $r \in \{G, N, B\}$. For example, under image scoring (Fig. 1(C)), $\Phi_{C,r} = 1 - \epsilon_a$ and $\Phi_{D,r} = \epsilon_a$ for any r . Under the so-called scoring-standing (see Sec. 4.1 for the notation of the social norms) shown in Fig. 1(C), $\Phi_{C,r} = 1 - \epsilon_a$ for any r , $\Phi_{D,+} = \Phi_{D,0} = \epsilon_a$, and $\Phi_{D,-} = 1 - \epsilon_a$. Each term on the right-hand side of Eq. (3) is a multiplication of three probabilities, i.e., (i) probability p_r^* that a recipient with reputation r is selected, (ii) probability ζ_r or $1 - \zeta_r$ that a donor implements C or D, respectively, and (iii) probability $\Phi_{C,r}$ or $\Phi_{D,r}$ that the observer assigns + to the donor.

We substitute p_G^* , p_N^* , and p_B^* obtained from Eq. (2) in the right-hand side of Eq. (3), which we denote by $f(\Phi^*)$. We obtain Φ^* by solving $x = f(x)$, $0 \leq x \leq 1$. Because $f(x)$ has a quadratic numerator and denominator in terms of x , equation $x = f(x)$ has at most three

solutions. Under any pair of action rule and social norm, to which we refer as action–norm pair in the following, $0 < \epsilon_a \leq f(x) \leq (1 - \epsilon_i)(1 - \epsilon_a) < 1$ holds true for any $0 \leq x \leq 1$. Therefore, the iteration scheme, in which we start with an initial x value ($0 \leq x \leq 1$) and repeatedly apply f , always converges. If the iteration starting from $x = 0$ and that starting from $x = 1$ converge to the same value, the solution specified by Φ^* and $\{p_G^*, p_N^*, p_B^*\}$ is unique. We confirmed that $x = f(x)$ has a unique solution for each of the $8 \times 64 = 512$ action–norm pairs.

Let $\Psi(\sigma, \{p_G, p_N, p_B\})$ be the probability that a donor with action rule σ cooperates with a recipient randomly chosen from a population according to reputation distribution $\{p_G, p_N, p_B\}$. We obtain

$$\Psi(\sigma, \{p_G, p_N, p_B\}) = p_G \zeta_G + p_N \zeta_N + p_B \zeta_B. \quad (4)$$

The average payoff per round to a resident player in the homogeneous population is given by

$$\pi = -c\Psi(\sigma, \{p_G^*, p_N^*, p_B^*\}) + b\Psi(\sigma, \{p_G^*, p_N^*, p_B^*\}). \quad (5)$$

To examine the stability of the homogeneous population, we consider an infinitesimally small fraction of mutants with action rule $\sigma' (\neq \sigma)$ that invades the homogeneous resident population. The equilibrium probability distribution of the mutant's reputation, denoted by $\{p_G'^*, p_N'^*, p_B'^*\}$, the equilibrium probability that a mutant receives +, denoted by Φ'^* , and payoff to a mutant player, denoted by π' , are given by

$$(p_G'^*, p_N'^*, p_B'^*) = \left(\frac{\Phi'^{*2}}{1 - \Phi'^* + \Phi'^{*2}}, \frac{\Phi'^*(1 - \Phi'^*)}{1 - \Phi'^* + \Phi'^{*2}}, \frac{(1 - \Phi'^*)^2}{1 - \Phi'^* + \Phi'^{*2}} \right), \quad (6)$$

$$\Phi'^* = \sum_{r \in \{G, N, B\}} p_r^* [\zeta_r' \Phi_{C,r} + (1 - \zeta_r') \Phi_{D,r}], \quad (7)$$

and

$$\pi' = -c\Psi(\sigma', \{p_G^*, p_N^*, p_B^*\}) + b\Psi(\sigma, \{p_G'^*, p_N'^*, p_B'^*\}), \quad (8)$$

respectively. In Eq. (7), ζ_r' represents the probability that a mutant cooperates with a recipient with reputation r . It should be noted that p_r^* in Eq. (7) is the solution of Eqs. (2) and (3) and that $\Phi_{C,r}$ and $\Phi_{D,r}$ also refer to the values for the resident population.

Action rule σ adopted by the resident players is strict Nash equilibrium if the payoff to a resident player (i.e., π) is larger than the payoff to any mutant player (i.e., π'). Using Eqs. (5)

and (8), we obtain this condition as follows:

$$\begin{aligned} & \frac{b}{c} [\Psi(\sigma, \{p_G^*, p_N^*, p_B^*\}) - \Psi(\sigma, \{p_G^*, p_N^*, p_B^*\})] \\ & < \Psi(\sigma', \{p_G^*, p_N^*, p_B^*\}) - \Psi(\sigma, \{p_G^*, p_N^*, p_B^*\}) \quad \text{for any } \sigma' \neq \sigma. \end{aligned} \quad (9)$$

We also investigate the stability of action rule σ against invasion by a previously identified strong competitor (Leimar & Hammerstein, 2001), which is so-called the Self strategy (Ohtsuki & Iwasa, 2004). A Self donor plays a donation game as selfishly as possible under the constraint that the donor's reputation stays above a threshold. When determining the action, the Self donor refers to its own reputation and does not refer to the recipient's reputation. We assume that the Self donor defects if its reputation is G and cooperates otherwise. Under the gradual reputation dynamics, the Self player maintains its reputation value at G or N, not B, except in the case of error.

3.2 Heterogeneous Populations Composed of Two Action Rules

We also examine the stability of heterogeneous populations in which two action rules, denoted by σ_1 and σ_2 , coexist with fractions q_1 and q_2 , respectively ($q_1 + q_2 = 1$). For each of $(q_1, q_2) = (0.01, 0.99), (0.02, 0.98), \dots, (0.99, 0.01)$, we first calculate the equilibrium probabilities of + for σ_1 and σ_2 by an iteration scheme similar to that described in Sec. 3.1. Under the gradual reputation dynamics, the distribution of the reputation values is given by

$$(p_{G,i}^*, p_{N,i}^*, p_{B,i}^*) = \left(\frac{\Phi_i^{*2}}{1 - \Phi_i^* + \Phi_i^{*2}}, \frac{\Phi_i^*(1 - \Phi_i^*)}{1 - \Phi_i^* + \Phi_i^{*2}}, \frac{(1 - \Phi_i^*)^2}{1 - \Phi_i^* + \Phi_i^{*2}} \right), \quad (10)$$

where $p_{r,i}^*$ is the equilibrium probability that a player with action rule i ($i = 1, 2$) possesses reputation $r \in \{G, N, B\}$, and Φ_i^* represents the equilibrium probability that a σ_i player receives +. The equivalent of Eq. (10) and the following results can be obtained similarly for the saltatory reputation dynamics. Φ_i^* is given by

$$\Phi_i^* = \sum_{j=1}^2 \sum_{r \in \{G, N, B\}} q_j p_{r,j}^* [\zeta_{r,i} \Phi_{C,r} + (1 - \zeta_{r,i}) \Phi_{D,r}], \quad (11)$$

where $\zeta_{r,i}$ is the probability that a σ_i donor implements C when the recipient has reputation r . Substitution of Eq. (10) into Eq. (11) leads to $\vec{\Phi}^* = g(\vec{\Phi}^*)$, where $\vec{\Phi}^* = (\Phi_1^*, \Phi_2^*)$.

It should be noted that $\vec{\Phi}^* = g(\vec{\Phi}^*)$ may have multiple fixed points. An example is given by the population composed of the equal fraction of $\sigma_1 = \text{gDisc}$ and $\sigma_2 = \text{rDisc}$, i.e., $q_1 = q_2 = 0.5$, under image scoring. In this case, both a cooperative population (i.e., $\Phi_1^*, \Phi_2^* \approx (1 - \epsilon_i)(1 - \epsilon_a)$) and a defective population (i.e., $\Phi_1^*, \Phi_2^* \approx \epsilon_a$) satisfy $\vec{\Phi}^* = g(\vec{\Phi}^*)$. Because of the multistability, we adopt $11^2 = 121$ initial conditions, i.e., $\vec{\Phi} = (0.1i, 0.1j)$, $0 \leq i, j \leq 10$, for the iteration scheme to identify all the fixed points. In fact, we find that the multistability does not cause a severe problem. We will show in Results that five mixed populations composed of two action rules are stable and realize a sufficiently large probability of cooperation. Among them, only one population, which consists of gDisc and rDisc and is stable under the so-called scoring-shunning social norm, yields the bistable equilibria. One equilibrium yields a large cooperation probability (> 0.93) and the other equilibrium yields a low cooperation probability (< 0.03). These results are qualitatively the same for the gradual and saltatory reputation dynamics. For this social norm, we only keep the more cooperative equilibrium.

Then, we obtain the trinary distribution of reputation for each of the two action rules by substituting Φ_i^* in Eq. (11). The average payoff per round to a σ_i resident player ($i = 1, 2$) is given by

$$\pi_i = -c \sum_{j=1}^2 q_j \Psi(\sigma_i, \{p_{G,j}^*, p_{N,j}^*, p_{B,j}^*\}) + b \sum_{j=1}^2 q_j \Psi(\sigma_j, \{p_{G,i}^*, p_{N,i}^*, p_{B,i}^*\}), \quad (12)$$

where $\Psi(\sigma_i, \{p_{G,j}^*, p_{N,j}^*, p_{B,j}^*\})$ is the probability that a σ_i donor cooperates with a σ_j recipient. In the equilibrium, the payoffs to the players with the different action rules are the same. Therefore, we calculate the value of b/c for which $\pi_1 = \pi_2$. If $d(\pi_1 - \pi_2)/dq_1 > 0$, the mixed population is unstable against an infinitesimally small drift of the fraction of the two action rules. We are concerned with the pairs of σ_1 and σ_2 that satisfy $d(\pi_1 - \pi_2)/dq_1 < 0$ at the obtained b/c value.

The mixed population is strict Nash when the payoff to any of the six mutants with a third action rule is smaller than that to a resident player. By substituting Eq. (6) in

$$\Phi'^* = \sum_{j=1}^2 \sum_{r \in \{G, N, B\}} q_j p_{r,j}^* [\zeta'_r \Phi_{C,r} + (1 - \zeta'_r) \Phi_{D,r}], \quad (13)$$

we obtain the equilibrium probability that a mutant receives + (i.e., Φ'^*) and the equilibrium

distribution of the mutant's reputation (i.e., $\{p_G^*, p_N^*, p_B^*\}$). The payoff to a mutant player is given by

$$\pi' = -c \sum_{j=1}^2 q_j \Psi(\sigma', \{p_{G,j}^*, p_{N,j}^*, p_{B,j}^*\}) + b \sum_{j=1}^2 q_j \Psi(\sigma_j, \{p_G^*, p_N^*, p_B^*\}). \quad (14)$$

As in the analysis of the homogeneous population, we also examine the stability of the heterogeneous population against invasion by the Self mutant (Sec. 3.1).

We do not examine the mixed population composed of more than two action rules.

4 Results

We refer to each action rule by concatenating three letters, either C or D. The first, second, and the third letters represent the intended action toward a recipient with reputation G, N, and B, respectively. For example, gDisc = CCD and rDisc = CDD.

4.1 Gradual Reputation Dynamics

4.1.1 Enumeration of Stable Populations

We set $\epsilon_i = \epsilon_a = 0.02$ and consider the gradual reputation dynamics (Fig. 2(A)). We found that the homogeneous population is stable against invasion by mutants for some benefit-to-cost values b/c for 108 out of $8 \times 64 = 512$ action–norm pairs. We exclude 64 pairs with $\sigma = \text{AllD}$ from the 108 pairs because of the lack of cooperation. The entire game is symmetric with respect to the simultaneous flipping of $G \leftrightarrow B$ and $+ \leftrightarrow -$ (Ohtsuki & Iwasa, 2004). Therefore, there are $(108 - 64)/2 = 22$ essentially distinct pairs. The 22 pairs are listed in Table 1. In addition, there are nine essentially distinct mixtures of two action rules that are stable under a certain social norm. Among these stable populations, there are 12 homogeneous populations composed of a single action rule and five heterogeneous populations composed of two action rules that realize a probability of cooperation larger than 0.5 for a b/c value smaller than 20. Our additional numerical simulations suggest that the probability of cooperation tends to unity if and only if this criterion is satisfied (Appendix B).

The cooperative equilibria, i.e., stable populations satisfying this criterion under a given social norm, are summarized in Fig. 3(A). In Fig. 3(A), a social norm is represented by a combination of s_G , s_N , and s_B , each of which takes either $++$, $-+$, $--$, or $+-$. For example, $s_G = -+$ indicates that donor's implemented action C and D toward a G recipient is assessed to be $-$ and $+$, respectively. In fact, all the cooperative equilibria require $s_G = +-$ such that only s_N and s_B are indicated in Fig. 3(A).

4.1.2 $s_N = ++$ or $-+$

Six action–norm pairs having $\sigma = \text{rDisc}$, $s_G = +-$, $s_N = ++$ or $-+$, and $s_B = ++$, $-+$, or $--$ realize a large probability of cooperation (≈ 0.94), a large probability of $+$ (≈ 0.98), and a mild restriction on b/c (i.e., $b/c > 1.004$). In these equilibria, most resident players have the G reputation because the probability of $+$ is large. A small fraction of players possesses the N reputation owing to error (see Sec. 2.1 for the definition of two types of error), and an even smaller fraction of players possesses the B reputation. We verified that the population is stable against invasion by Self players (Sec. 3.1) under each of the six social norms.

We call the six social norms standing–standing, standing–judging, standing–shunning, judging–standing, judging–judging, and judging–shunning. The first (second) half of the name represents the social norm represented by s_G and s_N (s_G and s_B) in the case of the binary reputation model. For example, the combination of $s_G = +-$ and $s_N = ++$ represents the standing social norm in the binary model if N is identified with B in the binary model (Fig. 1(D)). Similarly, the combination of $s_G = +-$ and $s_B = -+$ represents the judging social norm in the binary model. Therefore, we call the social norm given by $s_G = +-$, $s_N = ++$, and $s_B = -+$ standing–judging.

The six norms realize nearly perfect cooperation because the first half of the social norm (i.e., combination of s_G and s_N) is either standing or judging and the rDisc donor cooperates with G recipients, but not N recipients. It should be noted that standing and judging are the only second-order social norms that stabilize cooperation without special conditions in the binary reputation model (Ohtsuki & Iwasa, 2004). Even if the donor's reputation transits from N to B owing to error, the donor regains the N reputation by cooperating with a G recipient,

which occupies the majority of the population. Although donors meeting recipients with the B reputation receive $-$ under $s_B = --$, cooperation at the population level is maintained because few players have the B reputation.

4.1.3 $s_N = --$

Two cooperative action–norm pairs, i.e., $\sigma = \text{rDisc}$, $s_G = +-$, $s_N = --$, and $s_B = ++$ or $-+$, are also stable under a mild condition on b/c , i.e., $b/c > 1.018$. The homogeneous population of rDisc is not invaded by Self strategy under each of the two social norms. However, the probability of $+$ (≈ 0.74) and that of cooperation (≈ 0.66) are not as large as those for the previous six action–norm pairs. This is intuitively because, under the current social norms, i.e., shunning–standing and shunning–judging, a donor whose recipient has an N reputation always receives $-$ except in the case of the assessment error. This behavior of the model is similar to that under shunning in the binary reputation model. Nevertheless, different from the case of the binary model, the cooperation probability tends to unity in the error-free limit (Appendix B) for the two norms in the trinary reputation model. We discuss this point further in Discussion.

4.1.4 $s_N = +-$

The other cooperative equilibria are four homogeneous populations and five heterogeneous populations, which are stable under social norms satisfying $s_G = s_N = +-$ (Fig. 3(A)). For $s_B = ++$ (scoring–standing) and $-+$ (scoring–judging), a homogeneous population composed of gDisc is stable for large b/c (> 8.480), and a mixed population composed of gDisc and rDisc is stable for small b/c (< 8.480). For $s_B = --$ (scoring–shunning), a homogeneous population composed of gDisc is stable for large b/c (> 8.108), and a mixed population composed of gDisc and rDisc is stable for small b/c (< 8.108). The fraction of gDisc increases with b/c under scoring–standing and is unity when $b/c > 8.480$, as shown in Fig. 3(B). The results shown in Fig. 3(B) are indistinguishable from those for scoring–judging and similar to those for scoring–shunning. For the three social norms, the fraction of gDisc converges to $1 - c/b$ in the limit $\epsilon_i, \epsilon_a \rightarrow 0$ (Appendix C), which implies that only the mixed population of gDisc and rDisc

is stable in the error-free limit. Consistent with this, the threshold value of b/c above which the homogeneous gDisc population is stable diverges as the error probabilities become small (Appendix B).

If we hypothetically merge G and N, our result that cooperation is stable under scoring-standing and scoring-judging corresponds to the fact that the cooperation is stable under standing and judging, respectively, in the binary model (Ohtsuki & Iwasa, 2004). Under scoring-shunning, cooperation is not undermined by $s_B = --$ for the same reason as that for standing-shunning and judging-shunning (see Sec. 4.1.2 and Discussion).

Under scoring-scoring (i.e., $s_G = s_N = s_B = +-$), which we also call the image scoring (Fig. 1(C)), there are three types of stable populations depending on the b/c value. When $1.941 < b/c < 8.230$, the mixed population composed of gDisc and CDC is stable (Fig. 3(C)). We regard CDC as a variant of rDisc because there are few players having the B reputation in the stable population; a CDC player obtains a slightly larger payoff than an rDisc player. In fact, the mixed population of gDisc and rDisc is stable against invasion by all but CDC mutants. When $8.230 < b/c < 12.53$, the homogeneous population of gDisc is stable (Fig. 3(C)). When $b/c > 12.53$, the mixed population composed of gDisc and AllC is stable (Fig. 3(C)). In all the three cases, cooperation occurs with a large probability (> 0.94). This result is in a stark contrast with that in the binary reputation model, whereby image scoring does not usually support cooperation (Panchanathan & Boyd, 2003; Ohtsuki, 2004; Ohtsuki & Iwasa, 2004; Ohtsuki & Iwasa, 2007; Roberts, 2008). We discuss this discrepancy in Discussion.

Under image scoring, the values of b/c and the error probabilities under the assumption $\epsilon_i = \epsilon_a$, for which one of the three populations is stable are shown in Fig. 4(A). For small error, the homogeneous population of gDisc and the mixed population composed of gDisc and AllC are invaded by rDisc mutants. Only the heterogeneous population composed of gDisc and CDC survives the error-free limit. In this limit, the fraction of gDisc is given by $1 - c/b$ for the same reason as that for the heterogeneous population composed of gDisc and rDisc under scoring-standing (Sec. 4.1.4; also see Appendix C).

The nine populations stable under $s_G = s_N = +-$ are unstable against invasion by Self mutants. When there are more than two possible reputation values, Self is generally a strong com-

petitor under image scoring such that it undermines cooperation (Leimar & Hammerstein, 2001). Our results are consistent with theirs. In the next section, we propose a different reputation dynamics that makes cooperation stable against invasion by Self under image scoring.

4.2 Saltatory Reputation Dynamics

4.2.1 Downward Saltation

We first investigate a saltatory reputation dynamics (Fig. 2(B)) given by $\beta_d = 0.5$ and $\beta_u = 0$. This dynamics implies that a G reputation of a donor jumps down to a B reputation in one step with a probability $\beta_d = 0.5$. We set $\epsilon_i = \epsilon_a = 0.02$ and identify stable populations under each social norm. To select cooperative equilibria, we imposed the same condition as that in the case of the gradual reputation dynamics, i.e., a cooperation probability larger than 0.5 for some $b/c < 20$.

We find that the set of cooperative equilibria (without consideration of Self mutants) is the same as that for the gradual reputation dynamics (Fig. 3(A)). The cooperation probability is equal to 0.9237 for rDisc under standing–standing, 0.8614 for rDisc under standing–shunning, 0.6804 for rDisc under shunning–standing, 0.9340 for gDisc under scoring–standing, and 0.9340 for gDisc under scoring–shunning or image scoring. The cooperation probability is preserved if we replace standing by judging. Because a G reputation can turn into a B reputation in one step, the probability of cooperation is somewhat smaller than in the case of the gradual reputation dynamics. The difference in the cooperation probability between the two reputation dynamics is relatively large when s_B is $--$ or $+-$. This is because players with the B reputation, albeit occupying a small fraction, can make other players with the G reputation to transit to the B reputation in one step under these social norms.

Under image scoring, the stable population as a function of b/c and the error probabilities under the assumption $\epsilon_i = \epsilon_a$ is shown in Fig. 4(B). Figure 4(B) is similar to the results for the gradual reputation dynamics shown in Fig. 4(A). As a slight difference between Figs. 4(A) and 4(B), less cooperative populations (e.g., a homogeneous population of gDisc as compared to a heterogeneous population of gDisc and AllC) can be stable in the downward saltatory

reputation dynamics than in the gradual reputation dynamics for the same values of b/c and the error probability. This is intuitively because the downward saltatory reputation dynamics yields worse reputations than the gradual reputation dynamics.

We find that Self mutants cannot invade any of the stable populations shown in Fig. 3(A).

4.2.2 Upward Saltation

We next identified the stable populations under the saltatory reputation dynamics (Fig. 2(B)) given by $\beta_d = 0$, and $\beta_u = 0.5$. This dynamics implies that a B reputation of a donor jumps up to a G reputation in one step with a probability $\beta_u = 0.5$. We set $\epsilon_i = \epsilon_a = 0.02$ and identify stable populations under each social norm. To select cooperative equilibria, we imposed the same condition as that in the case of the gradual reputation dynamics, i.e., a cooperation probability larger than 0.5 for some $b/c < 20$.

We find that the set of cooperative equilibria (without consideration of Self mutants) is the same as that for the gradual reputation dynamics. The cooperation probabilities in stable and cooperative populations are equal to 0.9416 for rDisc under standing-standing, 0.9401 for rDisc under standing-shunning, 0.7605 for rDisc under shunning-standing, 0.9785 for gDisc under scoring-standing, and 0.9783 for gDisc under scoring-shunning or image scoring. The cooperation probability is preserved if we replace standing by judging. Because a B reputation can turn into a G reputation in one step, the probability of cooperation is a little larger than in the case of the gradual reputation dynamics. The difference is relatively large when $s_N = --$.

Under image scoring, the stable population as a function of b/c and ϵ_i ($= \epsilon_a$) is shown in Fig. 4(C). Figure 4(C) is quantitatively different from the results for the gradual reputation dynamics shown in Fig. 4(A). More cooperative action rules (e.g., a heterogeneous population of gDisc and AllC compared to a homogeneous population of gDisc) are stable in the upward saltatory reputation dynamics than in the gradual reputation dynamics for the same values of b/c and the error probability. This is intuitively because the upward saltatory reputation dynamics yields better reputations than the gradual reputation dynamics.

We find that Self mutants cannot invade a homogeneous population of rDisc under each of the eight social norms satisfying $s_G = +-$ and $s_N = ++, -+, --$. However, under social

norms satisfying $s_G = s_N = +-$, with image scoring included, the nine populations are not stable against invasion by Self mutants. The result that Self undermines cooperation is the same as that for the gradual reputation dynamics (Fig. 4(A)).

5 Discussion

We analyzed a trinary reputation model of indirect reciprocity and identified cooperative Nash equilibria composed of a single action rule or mixture of two action rules. Independent of details of the reputation dynamics (i.e., gradual or saltatory), we found at a small error level (i.e., $\epsilon_i = \epsilon_a = 0.02$) that 12 homogeneous populations and five mixed populations are cooperative and stable under different social norms (Fig. 3(A)). When the error probabilities are even smaller, eight homogeneous populations and four mixed populations remain stable (Appendix B). In particular, under image scoring, the heterogeneous population composed of gDisc and CDC (a variant of rDisc) is stable in the error-free limit (Fig. 4(A); see Fig. 1(B) for the definition of gDisc, CDC, and rDisc).

The results derived from the trinary reputation model and those derived from the binary reputation model (Nowak & Sigmund, 1998a) are similar in some aspects. For example, the standing-standing social norm (i.e., $s_G = +-$, $s_N = s_B = ++$) in the trinary model coincides with the standing social norm in the binary model (Fig. 1(D)) if we merge the N and B reputations in the trinary model. Therefore, the result that standing-standing enables cooperation is not surprising. Similarly, scoring-judging (i.e., $s_G = s_N = +-$, $s_B = -+$), for example, is almost equivalent to judging in the binary model if we merge G and N reputations in the trinary model. However, the results for the trinary and binary models are fundamentally different in the following two aspects.

First, shunning (see Fig. 1(D) in the case of the binary reputation model) is more supportive of cooperation in the trinary than binary model. In the binary model, shunning results in a cooperation probability of $\approx 1/2$ in the error-free limit unless a cooperation prone initial condition and a finite number of rounds are combined (Ohtsuki & Iwasa, 2007) or reputations of players are only partially visible to others (Nakamura & Masuda, 2011). In the trinary

model, shunning–standing and shunning–judging stabilize full cooperation in the error-free limit (Appendix B). Under these social norms, a donor meeting a recipient with an N reputation receives $-$ irrespective of the action, which is common to the behavior of the binary model. In the trinary model, however, even if a donor obtains a B reputation, the donor easily receives $+$ either by cooperating with a G recipient or justifiably defecting against a B recipient. Owing to the contribution of the justified defection, a donor gains $+$ more often than $-$ on average. If players with the G reputation increase in number owing to this mechanism, players more likely receive $+$ than $-$. This positive feedback sustains cooperation under shunning–related social norms in the trinary model.

Second, and the more important, image scoring is capable of supporting cooperation in our model. Even if we consider Self mutants, which avoid a B reputation and are as selfish as possible, image scoring supports cooperation if the probability that a G reputation transits to a B reputation in one step is positive. In previous literature, the Self strategy is recognized as a strong competitor that spoils cooperation in the indirect reciprocity game with more than two possible reputation values (Leimar & Hammerstein, 2001). Our conclusion that image scoring can support cooperation is consistent with the results derived from behavioral experiments (Wedekind & Milinski, 2000; Milinski *et al.*, 2001; Seinen & Schram, 2006) and those derived from numerical simulations (Nowak & Sigmund, 1998b; Diekmann & Przepiorka, 2005) but opposite to those derived from the binary reputation model (Panchanathan & Boyd, 2003; Ohtsuki, 2004; Ohtsuki & Iwasa, 2004; Ohtsuki & Iwasa, 2007). We reached this conclusion by simply introducing a third reputation to the standard binary model of indirect reciprocity. The cooperation under image scoring in our model does not require forced cooperation in the first round (Nowak & Sigmund, 1998a), partial cooperation of defective players (Fishman, 2003), binomially or Poisson distributed number of rounds (Brandt & Sigmund, 2004; Brandt & Sigmund, 2006), growth of social networks used for transmission of reputation (Brandt & Sigmund, 2005), or a small probability with which the donor’s reputation is revealed to other players (Uchida, 2010).

In our model, cooperation under image scoring occurs for the following intuitive reason. Although the composition of the stable population depends on the benefit-to-cost ratio and the error probabilities, the main action rule present in the stable population is gDisc, which

cooperates with G and N recipients and defects against B recipients. In the equilibrium, most gDisc resident players possess the G reputation. If some AllD mutants are present, a B player in the population is likely to be AllD and not gDisc. With the ternary reputation, gDisc players justifiably defect (i.e., D against a B recipient) but do not selfishly defect (i.e., D against a G or N recipient). A donor that has justifiably defected would receive an N reputation but not a B reputation because few recipients have the B reputation in the population. Therefore, gDisc players would not obtain the B reputation. gDisc players that happen to obtain the N reputation, when selected as donor, likely meet a recipient with a G reputation such that they regain the G reputation. In contrast, AllD players would obtain the B reputation such that they are not helped by others. This contrasts with the case of the binary reputation model, in which a discriminating donor that defects against whatever recipient immediately receives the worst reputation (i.e., B) and then is defected by others.

The presence of downward saltation (i.e., transition from the G reputation to B reputation in one step) is a key to make a cooperative population stable against invasion by Self mutants. To explain this point intuitively, let us consider a homogeneous population of gDisc and assume that the implementation and assessment errors are absent. If downward saltation is absent in the reputation dynamics, Self mutants flip between G and N reputations by alternatingly cooperating and defecting. By doing so, the Self mutants can elicit cooperation from gDisc residents and cooperate with probability $\approx 1/2$. Because gDisc residents cooperate with probability ≈ 1 , the gDisc population is invaded by Self mutants. However, if downward saltation can occur in the reputation dynamics, a nonnegligible fraction of Self mutants possess the B reputation because they defect when they have the G reputation. In contrast, gDisc players maintain the G reputation because they cooperate even if they have the G reputation. Therefore, in the presence of downward saltation, Self mutants cannot invade the population of gDisc residents.

The mechanism of cooperation under image scoring in our model is similar to that under the so-called tolerant scoring proposed by Berger (Berger, 2011). The population of tolerant discriminator, which defects against a recipient if the recipient has defected in the last two rounds and cooperates otherwise, is stable. Berger’s results and ours are different in the following aspects. First, Berger assumed three action rules, i.e., tolerant discriminator, which is similar

to gDisc, AllC, and AllD and investigated the behavior of the model under tolerant scoring. We carried out an exhaustive search of the space of the action rule and social norms to find the viability of image scoring. Second, in our model, cooperation is stable even if the reputation moves only one step in a round, if the Self mutants are not considered. It should be noted that Berger did not consider the Self mutants. Third, cooperation is realized by a homogeneous population of one type of players in the Berger's model. In our model, cooperation is realized by mixture of two types of discriminative players (i.e., gDisc and CDC).

The present study has following limitations. First, we assumed that all the players in the population use the same social norm. This oversimplification excludes a possible situation in which different norms compete in a population (e.g., (Pacheco *et al.*, 2006; Uchida, 2010)). Second, we only analyzed the stability of populations composed of up to two action rules for simplicity. Third, we analyzed local stability of the equilibria and disregarded dynamics. Even when a cooperative equilibrium is locally stable, it may in fact attract a tiny fraction of initial conditions. Fourth, we assumed that the reputation sharing is public. In other words, the information about the donor's new reputation immediately spreads from the observer to the entire population. In a large population, such immediate spreading is impossible, and one has to assume, for example, that only a fraction of players gains the information about a donor in one game (Nowak & Sigmund, 1998b; Brandt & Sigmund, 2004; Uchida, 2010).

Appendix A: Distribution of the Reputation for the Saltatory Reputation Dynamics

Under the saltatory reputation dynamics (Fig. 2(B)), we obtain

$$\begin{cases} p_G^* &= p_G^* \Phi^* + p_N^* \Phi^* + p_B^* \beta_u \Phi^*, \\ p_N^* &= p_G^* (1 - \beta_d) (1 - \Phi^*) + p_B^* (1 - \beta_u) \Phi^*, \\ p_B^* &= p_G^* \beta_d (1 - \Phi^*) + p_N^* (1 - \Phi^*) + p_B^* (1 - \Phi^*). \end{cases} \quad (15)$$

Equation (15) and the normalization $p_G^* + p_N^* + p_B^* = 1$ lead to

$$\begin{bmatrix} p_G^* \\ p_N^* \\ p_B^* \end{bmatrix} = \frac{1}{Z} \begin{bmatrix} (1 - \beta_u)\Phi^{*2} + \beta_u\Phi^* \\ (1 - \beta_d\beta_u)\Phi^*(1 - \Phi^*) \\ (1 - \beta_d)(1 - \Phi^*)^2 + \beta_d(1 - \Phi^*) \end{bmatrix}, \quad (16)$$

where $Z = 1 - (1 - \beta_d)(1 - \beta_u)\Phi^*(1 - \Phi^*)$.

Appendix B: Justification of the Criterion of Equilibrium Selection

We defined cooperative equilibrium as stable population in which the probability of cooperation is larger than 0.5 for some b/c such that $b/c < 20$. To justify this criterion, we carry out additional numerical simulations with various error probabilities satisfying $\epsilon_i = \epsilon_a$. The probability of cooperation in each stable homogeneous population is shown for various error probability values in Fig. 5. Figure 5 shows that, under 13 out of the 22 stable action–norm pairs, the probability of cooperation seems to converge to unity in the error-free limit. It should be noted that we identified 12, not 13, cooperative homogeneous populations in the main text. One of the 13 action–norm pairs is excluded because it is stable only for large b/c values (i.e., $b/c > 52.63$) at $\epsilon_i = \epsilon_a = 0.02$ (Table 1). The probability of cooperation is larger than 0.5 for some b/c for five out of the nine stable heterogeneous populations composed of two action rules. For all the five heterogeneous populations, the probability of cooperation is larger than 0.5 for a b/c value smaller than 20 at $\epsilon_i = \epsilon_a = 0.02$.

In the homogeneous populations, under eight out of the 13 social norms for which the probability of cooperation seems to converge to unity in the error-free limit, i.e., standing–standing, standing–judging, standing–shunning, judging–standing, judging–judging, judging–shunning, shunning–standing, and shunning–judging, the range of b/c in which the corresponding action rule is stable tends to $b/c > 1$ as the error probabilities become small. For example, under standing–standing, standing–judging, judging–standing, and judging–judging, the homogeneous population of rDisc is stable in the range $b/c > 1.00006$ when $\epsilon_i = \epsilon_a = 0.0025$; the corresponding range is $b/c > 1.004$ when $\epsilon_i = \epsilon_a = 0.02$ (Table 1). Under the other five

social norms, i.e., standing–scoring, scoring–standing, scoring–judging, scoring–shunning, and scoring–scoring, stable cooperation requires a large value of b/c when the error probabilities are small. For example, under standing–scoring, the homogeneous population of CDC is stable in the range $b/c > 402.1$ when $\epsilon_i = \epsilon_a = 0.0025$; the corresponding range is $b/c > 52.63$ when $\epsilon_i = \epsilon_a = 0.02$ (Table 1). Under image scoring, the homogeneous population of gDisc is stable in the range $b/c > 66.51$ when $\epsilon_i = \epsilon_a = 0.0025$; the corresponding range is $8.230 < b/c < 12.53$ when $\epsilon_i = \epsilon_a = 0.02$ (Table 1, Fig. 4(A)).

In the heterogeneous populations, four out of five equilibria, i.e., mixture of gDisc and rDisc under scoring–standing, mixture of gDisc and rDisc under scoring–judging, mixture of gDisc and rDisc under scoring–shunning, and mixture of gDisc and CDC under image scoring, realize a large probability of cooperation in a wide range of b/c when $\epsilon_i = \epsilon_a = 0.02$ (Figs. 3(B) and 3(C)) and also do so when the error probabilities are smaller. For example, the heterogeneous population composed of gDisc and CDC is stable under image scoring in the range $1.994 < b/c < 66.51$ when $\epsilon_i = \epsilon_a = 0.0025$; the corresponding range is $1.941 < b/c < 8.230$ when $\epsilon_i = \epsilon_a = 0.02$ (Fig. 4(A)).

Appendix C: Fraction of Two Action Rules in the Error-free Limit

We consider a social norm given by $s_G = s_N = +-$ and $s_B = ++, -+,$ or $--$. In the limit $\epsilon_i, \epsilon_a \rightarrow 0$, the fraction of gDisc, denoted by q_{gDisc} , and that of rDisc converge to $1 - c/b$ and c/b , respectively, for the following reason. In the equilibrium, few players possess the N or B reputation. The behavior of rDisc players and that of gDisc players differ only toward recipients with reputation N. An rDisc donor with a G reputation defects against a recipient with an N reputation, receives $-$, and transits to the N reputation. If this rDisc player is selected as recipient, an rDisc donor defects and a gDisc donor cooperates. If selected as donor, this rDisc player would cooperate because most recipients have the G reputation, receive $+$, and transit to G. Because the rDisc player with an N reputation is selected as recipient once on average before selected as donor, the expected payoff to the rDisc donor during this period is equal to

$q_{\text{gDisc}}b$ on average. During the same period, a gDisc donor with a G reputation cooperates with a recipient with an N reputation, pays c , receives $+$, keeps a G reputation, and gains benefit b when selected as recipient. By equating the payoffs to rDisc and gDisc players, we obtain $q_{\text{gDisc}}b = -c + b$, which leads to $q_{\text{gDisc}} = 1 - c/b$.

Acknowledgements

We thank Mitsuhiro Nakamura for critical reading of the manuscript and the guidance on the literature on image scoring. N.M. acknowledges the support provided through Grants-in-Aid for Scientific Research (No. 23681033, and Innovative Areas “Systems Molecular Ethology” (No. 20115009)) from MEXT, Japan. This research is also supported by the Aihara Project, the FIRST program from JSPS, initiated by CSTP.

References

- Axelrod, R. 1984. *The Evolution of Cooperation*. Basic Books, NY.
- Berger, U. 2011. Learning to cooperate via indirect reciprocity. *Games Econ. Behav.* 72, 30–37.
- Bolton, G. E. & Katok, E. 2004. How effective are electronic reputation mechanisms? An experimental investigation. *Manage. Sci.* 50, 1587–1602.
- Bolton, G. E., Katok, E. & Ockenfels, A. 2005. Cooperation among strangers with limited information about reputation. *J. Public Econ.* 89, 1457–1468.
- Brandt, H. & Sigmund, K. 2004. The logic of reprobation: assessment and action rules for indirect reciprocation. *J. Theor. Biol.* 231, 475–486.
- Brandt, H. & Sigmund, K. 2005. Indirect reciprocity, image scoring, and moral hazard. *Proc. Natl. Acad. Sci. USA*, 102, 2666–2670.
- Brandt, H. & Sigmund, K. 2006. The good, the bad and the discriminator—Errors in direct and indirect reciprocity. *J. Theor. Biol.* 239, 183–194.

- Diekmann, A. & Przepiorka, W. 2005. The evolution of trust and reputation: Results from simulation experiments. In *Third ESSA Conference* pp. 1–7.
- Engelmann, D. & Fischbacher, U. 2009. Indirect reciprocity and strategic reputation building in an experimental helping game. *Games Econ. Behav.* 67, 399–407.
- Fishman, M. A. 2003. Indirect reciprocity among imperfect individuals. *J. Theor. Biol.* 225, 285–292.
- Kandori, M. 1992. Social norms and community enforcement. *Rev. Econ. Stud.* 59, 63–80.
- Keser, C. 2003. Experimental games for the design of reputation management systems. *IBM Syst. J.* 42, 498–506.
- Leimar, O. & Hammerstein, P. 2001. Evolution of cooperation through indirect reciprocity. *Proc. R. Soc. B*, 268, 745–753.
- Milinski, M., Semmann, D., Bakker, T. C. M. & Krambeck, H.-J. 2001. Cooperation through indirect reciprocity: image scoring or standing strategy? *Proc. R. Soc. B*, 268, 2495–2501.
- Milinski, M., Semmann, D. & Krambeck, H.-J. 2002. Reputation helps solve the ‘tragedy of the commons’. *Nature*, 415, 424–426.
- Mohtashemi, M. & Mui, L. 2003. Evolution of indirect reciprocity by social information: the role of trust and reputation in evolution of altruism. *J. Theor. Biol.* 223, 523–531.
- Nakamura, M. & Masuda, N. 2011. Indirect reciprocity under incomplete observation. *PLoS Comput. Biol.* 7, e1002113.
- Nowak, M. A. & Sigmund, K. 1998a. The dynamics of indirect reciprocity. *J. Theor. Biol.* 194, 561–574.
- Nowak, M. A. & Sigmund, K. 1998b. Evolution of indirect reciprocity by image scoring. *Nature*, 393, 573–577.
- Nowak, M. A. & Sigmund, K. 2005. Evolution of indirect reciprocity. *Nature*, 437, 1291–1298.

- Ohtsuki, H. 2004. Reactive strategies in indirect reciprocity. *J. Theor. Biol.* 227, 299–314.
- Ohtsuki, H. & Iwasa, Y. 2004. How should we define goodness?—reputation dynamics in indirect reciprocity. *J. Theor. Biol.* 231, 107–120.
- Ohtsuki, H. & Iwasa, Y. 2006. The leading eight: Social norms that can maintain cooperation by indirect reciprocity. *J. Theor. Biol.* 239, 435–444.
- Ohtsuki, H. & Iwasa, Y. 2007. Global analyses of evolutionary dynamics and exhaustive search for social norms that maintain cooperation by reputation. *J. Theor. Biol.* 244, 518–531.
- Pacheco, J. M., Santos, F. C. & Chalub, F. A. C. C. 2006. Stern-judging: A simple, successful norm which promotes cooperation under indirect reciprocity. *PLoS Comput. Biol.* 2, e178.
- Panchanathan, K. & Boyd, R. 2003. A tale of two defectors: the importance of standing for evolution of indirect reciprocity. *J. Theor. Biol.* 224, 115–126.
- Resnick, P. & Zeckhauser, R. 2002. Trust among strangers in internet transactions: Empirical analysis of eBay’s reputation system. *Adv. Appl. Microeconomics*, 11, 127–157.
- Resnick, P., Zeckhauser, R., Swanson, J. & Lockwood, K. 2006. The value of reputation on eBay: A controlled experiment. *Exp. Econ.* 9, 79–101.
- Roberts, G. 2008. Evolution of direct and indirect reciprocity. *Proc. R. Soc. B*, 275, 173–179.
- Seinen, I. & Schram, A. 2006. Social status and group norms: Indirect reciprocity in a repeated helping experiment. *Eur. Econ. Rev.* 50, 581–602.
- Sigmund, K. 2012. Moral assessment in indirect reciprocity. *J. Theor. Biol.* 299, 25–30.
- Trivers, R. L. 1971. The evolution of reciprocal altruism. *Q. Rev. Biol.* 46, 35–57.
- Uchida, S. 2010. Effect of private information on indirect reciprocity. *Phys. Rev. E*, 82, 036111.
- Wedekind, C. & Braithwaite, V. A. 2002. The long-term benefits of human generosity in indirect reciprocity. *Curr. Biol.* 12, 1012–1015.

Wedekind, C. & Milinski, M. 2000. Cooperation through image scoring in humans. *Science*, 288, 850–852.

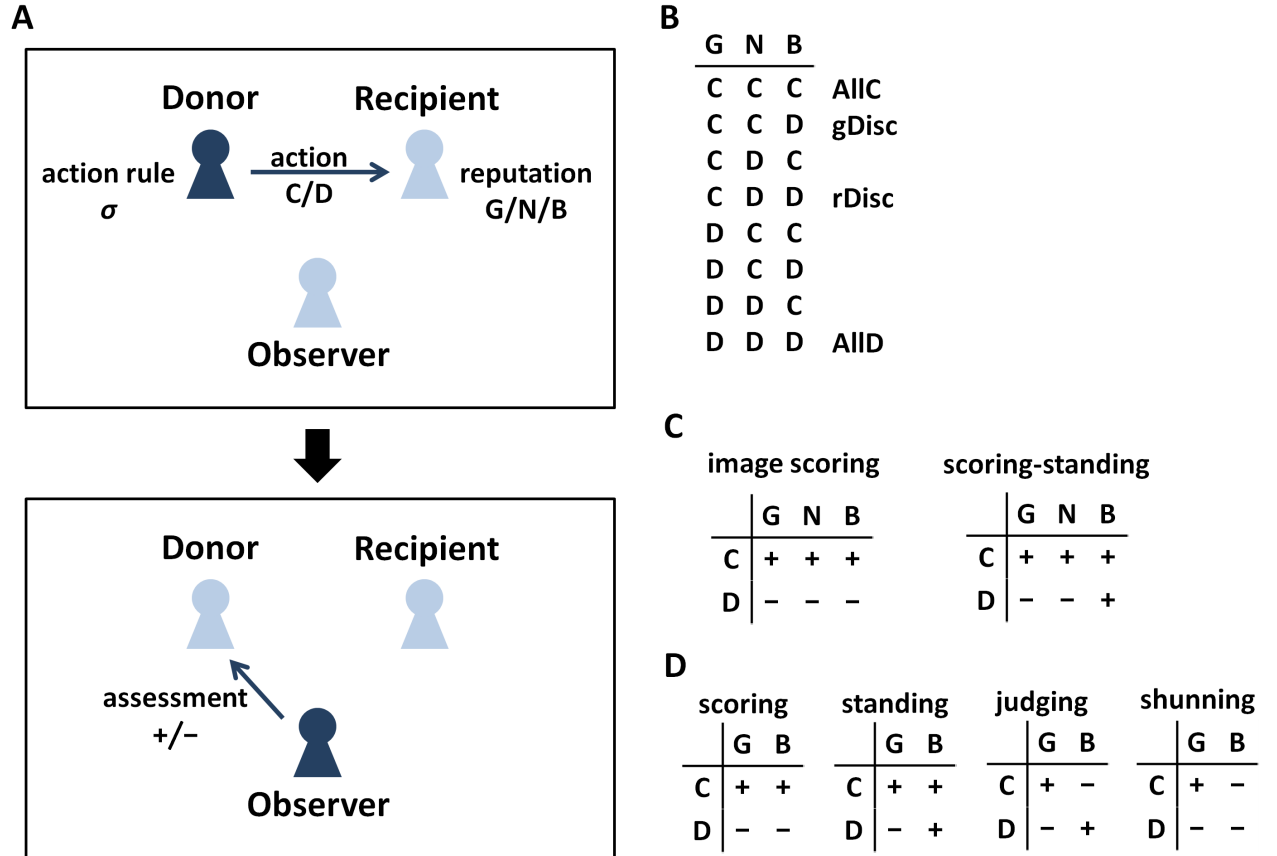


Figure 1: Rule of the donation game with trinary reputations. (A) Illustration of the interaction in a single game. (B) Eight action rules. (C) Representative social norms. The rows represent the donor's actions (i.e., C and D), the columns represent the recipient's reputations (G, N, and B), and + and - represent the assessments that observer assigns to the donor. (D) Representative social norms in the binary reputation model.

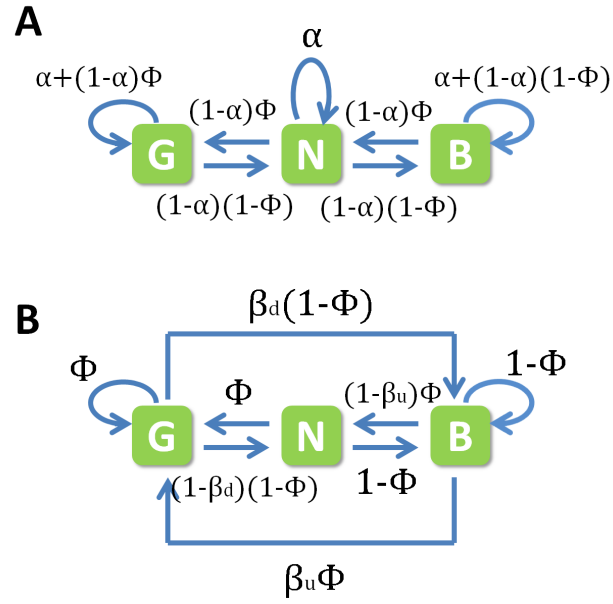


Figure 2: Two types of reputation dynamics. (A) Gradual reputation dynamics. (B) Saltatory reputation dynamics. Φ represents the probability that the donor receives +.

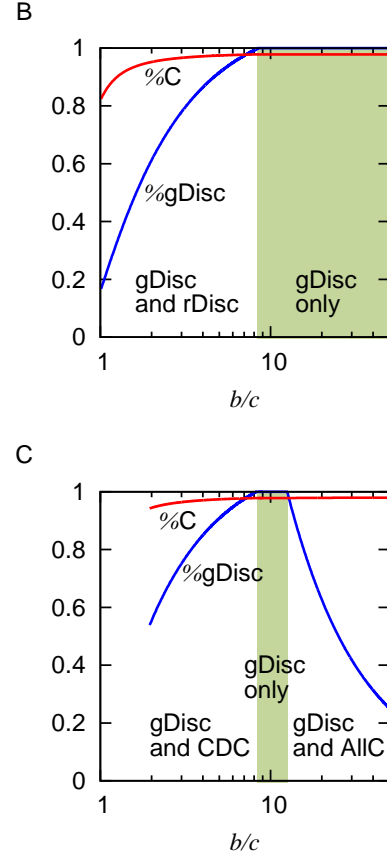
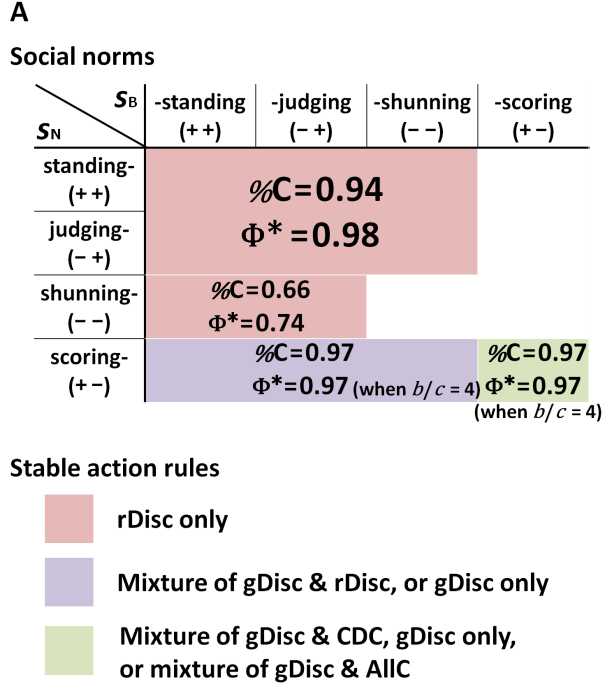


Figure 3: Results for the gradual reputation dynamics. We set $\epsilon_i = \epsilon_a = 0.02$. (A) Cooperative and stable action rules under different social norms. All the shown social norms own $s_G = +-$. (B) Average cooperation probability and the fraction of gDisc players under scoring-standing (i.e., $s_G = s_N = +-$ and $s_B = ++$). It should be noted that the fraction of gDisc and that of rDisc sum to unity. (C) Average cooperation probability and the fraction of gDisc players under scoring-scoring (also called image scoring; $s_G = s_N = s_B = +-$). The fraction of gDisc and that of CDC or AllC sum to unity.

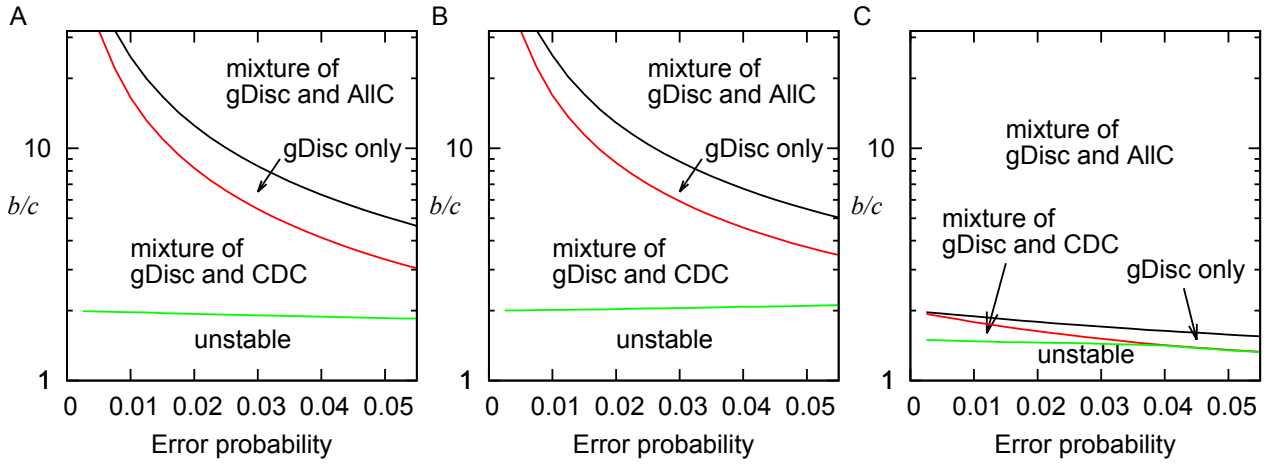


Figure 4: Parameter regions in which one of the three populations is stable under image scoring. We set the error probability $\epsilon_i = \epsilon_a$. (A) Gradual reputation dynamics, i.e., $(\beta_d, \beta_u) = (0, 0)$. (B) Saltatory reputation dynamics with downward saltation, i.e., $(\beta_d, \beta_u) = (0.5, 0)$. (C) Saltatory reputation dynamics with upward saltation, i.e., $(\beta_d, \beta_u) = (0, 0.5)$. The cooperative population is stable against invasion by Self mutants in the parameter region above the green line in (B). Self invades the cooperative population in all the parameter regions in (A) and (C).

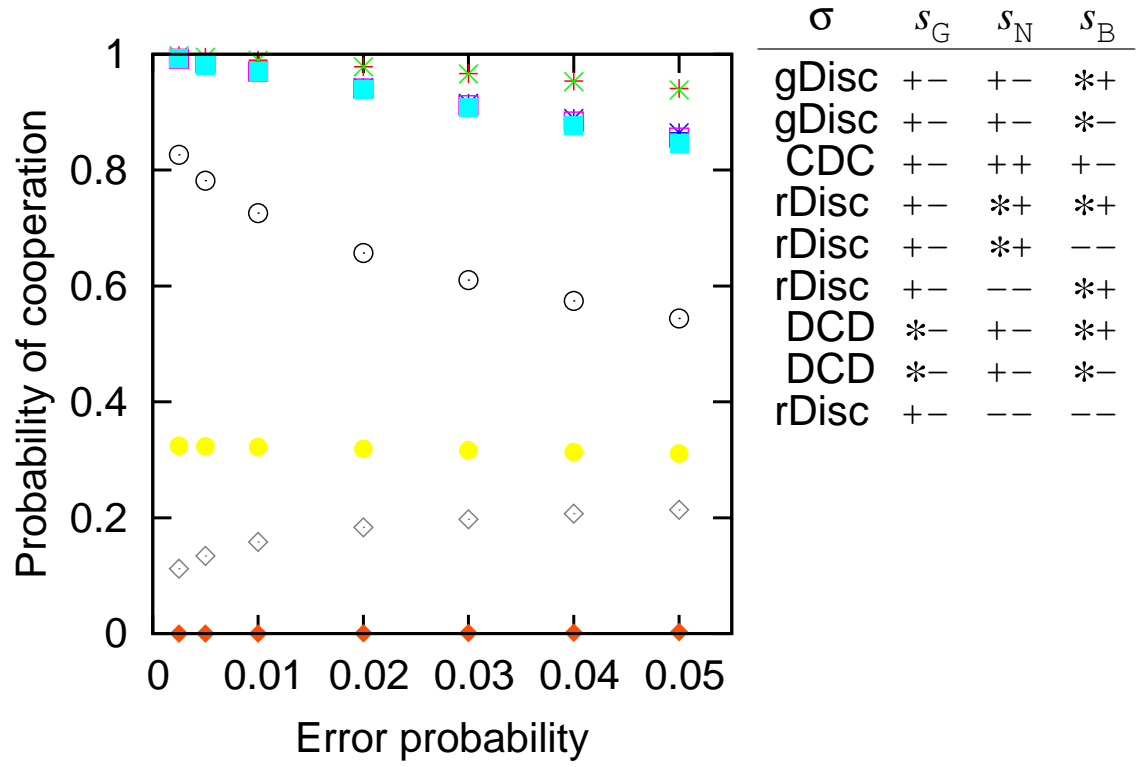


Figure 5: Relationships between the probability of cooperation and the error probability for the 22 stable homogeneous populations. We assume the gradual reputation dynamics and set $\epsilon_i = \epsilon_a$. An asterisk represents either $+$ or $-$.

Table 1: Stable action–norm pairs. We also show the range of b/c in which action rule σ is stable, probability of C, and mean + assessment (i.e., Φ^*) under the gradual reputation dynamics. An asterisk represents either + or –. We set $\epsilon_i = \epsilon_a = 0.02$.

| σ | s_G | s_N | s_B | range of b/c | %C | Φ^* |
|----------|-------|-------|-------|-----------------------|--------|----------|
| CCD | +- | +- | *+ | $8.480 < b/c$ | 0.9784 | 0.9800 |
| | | | -- | $8.108 < b/c$ | 0.9783 | 0.9783 |
| | | | +- | $8.230 < b/c < 12.53$ | | |
| CDC | +- | ++ | +- | $52.63 < b/c$ | 0.9424 | 0.9808 |
| CDD | +- | *+ | *+ | $1.004 < b/c$ | 0.9409 | 0.9808 |
| | | *+ | -- | $1.004 < b/c$ | 0.9392 | 0.9792 |
| | | -- | *+ | $1.018 < b/c$ | 0.6571 | 0.7445 |
| DCD | *- | +- | *+ | $3.716 < b/c$ | 0.3191 | 0.5688 |
| | +- | | *- | $1.137 < b/c < 1.296$ | 0.1833 | 0.1833 |
| | -- | | -- | $1.120 < b/c$ | | |
| | -- | | +- | $1.120 < b/c < 11.66$ | | |
| CDD | +- | -- | -- | $25.54 < b/c$ | 0.0004 | 0.0004 |